

Benchmarking saliency methods for chest X-ray interpretation

Received: 11 October 2021

Accepted: 26 August 2022

Published online: 10 October 2022

 Check for updates

Adriel Saporta^{1,9}, Xiaotong Gui^{2,9}, Ashwin Agrawal^{2,9}, Anuj Pareek³, Steven Q. H. Truong⁴, Chanh D. T. Nguyen^{4,5}, Van-Doan Ngo⁶, Jayne Seekins⁷, Francis G. Blankenberg⁷, Andrew Y. Ng², Matthew P. Lungren³ and Pranav Rajpurkar⁸  

Saliency methods, which produce heat maps that highlight the areas of the medical image that influence model prediction, are often presented to clinicians as an aid in diagnostic decision-making. However, rigorous investigation of the accuracy and reliability of these strategies is necessary before they are integrated into the clinical setting. In this work, we quantitatively evaluate seven saliency methods, including Grad-CAM, across multiple neural network architectures using two evaluation metrics. We establish the first human benchmark for chest X-ray segmentation in a multilabel classification set-up, and examine under what clinical conditions saliency maps might be more prone to failure in localizing important pathologies compared with a human expert benchmark. We find that (1) while Grad-CAM generally localized pathologies better than the other evaluated saliency methods, all seven performed significantly worse compared with the human benchmark, (2) the gap in localization performance between Grad-CAM and the human benchmark was largest for pathologies that were smaller in size and had shapes that were more complex, and (3) model confidence was positively correlated with Grad-CAM localization performance. Our work demonstrates that several important limitations of saliency methods must be addressed before we can rely on them for deep learning explainability in medical imaging.

Deep learning has enabled automated medical imaging interpretation at the level of practicing experts in some settings^{1–3}. While the potential benefits of automated diagnostic models are numerous, lack of model interpretability in the use of ‘black-box’ deep neural networks (DNNs) represents a major barrier to clinical trust and adoption^{4–6}. In fact, it has been argued that the European Union’s recently adopted General Data Protection Regulation affirms an individual’s right to an explanation in the context of automated decision-making⁷. Although the importance

of DNN interpretability is widely acknowledged and many techniques have been proposed, little emphasis has been placed on how best to quantitatively evaluate these explainability methods⁸.

One type of DNN interpretation strategy widely used in the context of medical imaging is based on saliency (or pixel-attribution) methods^{9–12}. Saliency methods produce heat maps highlighting the areas of the medical image that most influenced the DNN’s prediction. Since saliency methods provide post-hoc interpretability of models that are

¹Department of Computer Science, New York University, New York, NY, USA. ²Department of Computer Science, Stanford University, Stanford, CA, USA.

³Stanford Center for Artificial Intelligence in Medicine and Imaging, Palo Alto, CA, USA. ⁴VinBrain, Ha Noi, Vietnam. ⁵VinUniversity, Ha Noi, Vietnam.

⁶Vinmec International Hospital, Ha Noi, Vietnam. ⁷Department of Radiology, Stanford University School of Medicine, Stanford, CA, USA. ⁸Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA. ⁹These authors contributed equally: Adriel Saporta, Xiaotong Gui, Ashwin Agrawal.

 e-mail: pranav_rajpurkar@hms.harvard.edu

never exposed to bounding-box annotations or pixel-level segmentations during training, they are particularly useful in the context of medical imaging where ground-truth segmentations can be especially time-consuming and expensive to obtain. The heat maps help to visualize whether a DNN is concentrating on the same regions of a medical image that a human expert would focus on, rather than concentrating on a clinically irrelevant part of the medical image or even on confounders in the image^{13–15}. Saliency methods have been widely used for a variety of medical imaging tasks and modalities including, but not limited to, visualizing the performance of a convolutional neural network (CNN) in predicting (1) myocardial infarction¹⁶ and hypoglycemia¹⁷ from electrocardiograms, (2) visual impairment¹⁸, refractive error¹⁹ and anaemia²⁰ from retinal photographs, (3) long-term mortality²¹ and tuberculosis²² from chest X-ray (CXR) images and (4) appendicitis²³ and pulmonary embolism²⁴ on computed tomography scans. However, recent work has shown that saliency methods used to validate model predictions can be misleading in some cases and may lead to increased bias and loss of user trust in high-stakes contexts such as healthcare^{25–28}. Therefore, a rigorous investigation of the accuracy and reliability of these strategies is necessary before they are integrated into the clinical setting²⁹.

In this work, we perform a systematic evaluation of seven common saliency methods in medical imaging (Grad-CAM³⁰, Grad-CAM++³¹, ‘integrated gradients’³², Eigen-CAM³³, DeepLIFT³⁴, layer-wise relevance propagation (LRP)³⁵ and ‘occlusion’³⁶) using three common CNN architectures (DenseNet121³⁷, ResNet152³⁸ and Inception-v4³⁹). In doing so, we establish the first human benchmark for CXR segmentation in a multilabel classification set-up by collecting radiologist segmentations for ten pathologies using CheXpert, a large publicly available CXR dataset⁴⁰. To compare saliency method segmentations with expert segmentations, we use two metrics to capture localization accuracy: (1) mean intersection over union (mIoU), a metric that measures the overlap between the saliency method segmentation and the expert segmentation, and (2) hit rate, a less strict metric than mIoU that does not require the saliency method to locate the full extent of a pathology. We find that (1) while Grad-CAM generally localizes pathologies more accurately than the other evaluated saliency methods, all seven perform significantly worse compared with a human radiologist benchmark (although it is difficult to know whether poor localization performance is attributable to the model or to the saliency method), (2) the gap in localization performance between Grad-CAM and the human benchmark is largest for pathologies that are smaller in size and have shapes that are more complex, and (3) model confidence is positively correlated with Grad-CAM localization performance. We publicly release a development dataset of expert segmentations, which we call CheXlocalize, to facilitate further research in DNN explainability for medical imaging.

Results

Framework for evaluating saliency methods

Seven methods were evaluated—Grad-CAM, Grad-CAM++, integrated gradients, Eigen-CAM, DeepLIFT, LRP and occlusion—in a multilabel classification set-up on the CheXpert dataset (Fig. 1a). We ran experiments using three CNN architectures previously used on CheXpert: DenseNet121, ResNet152 and Inception-v4. For each combination of saliency method and model architecture, we trained and evaluated an ensemble of 30 CNNs (see Methods for ensembling details). We then passed each of the CXRs in the dataset’s holdout test set into the trained ensemble model to obtain image-level predictions for the following ten pathologies: ‘airspace opacity’, ‘atelectasis’, ‘cardiomegaly’, ‘consolidation’, ‘edema’, ‘enlarged cardiomediastinum’, ‘lung lesion’, ‘pleural effusion’, ‘pneumothorax’ and ‘support devices’. Of the 14 observations labelled in the CheXpert dataset, ‘fracture’ and ‘pleural other’ were not included in our analysis because they had low prevalence in our test set (fewer than ten examples), ‘pneumonia’ was not included because

it is a clinical (as opposed to a radiological) diagnosis and ‘no finding’ was not included because it is not applicable to evaluating localization performance. For each CXR, we used the saliency method to generate heat maps, one for each of the ten pathologies, and then applied a threshold to each heat map to produce binary segmentations (top row, Fig. 1a). Thresholding is determined per pathology using Otsu’s method⁴¹, which iteratively searches for a threshold value that maximizes interclass pixel intensity variance. We also conducted a second thresholding scheme in which we iteratively search for a threshold value that maximizes per-pathology mIoU on the validation set. There are no statistically significant differences between the two thresholding schemes when compared against the human benchmark (Extended Data Fig. 1). Additionally, to calculate the hit rate evaluation metric (described below), we extracted the pixel in the saliency method heat map with the largest value as the single most representative point on the CXR for that pathology.

We obtained two independent sets of pixel-level CXR segmentations on the holdout test set: ground-truth segmentations drawn by two board-certified radiologists (middle row, Fig. 1a) and human benchmark segmentations drawn by a separate group of three board-certified radiologists (bottom row, Fig. 1a). The human benchmark segmentations and the saliency method segmentations were compared with the ground-truth segmentations to establish the human benchmark localization performance and the saliency method localization performance, respectively. Additionally, for the hit rate evaluation metric, the radiologists who drew the benchmark segmentations were also asked to locate a single point on the CXR that was most representative of the pathology at hand (see Supplementary Figs. 1–11 for detailed instructions given to the radiologists). Note that the human benchmark localization performance demonstrates interrater variability, and we use it as a reference when evaluating saliency method pipelines.

We used two evaluation metrics to compare segmentations (Fig. 1b). Our primary metric, mIoU, measures how much, on average, either the saliency method or benchmark segmentations overlapped with the ground-truth segmentations. Our secondary metric, hit rate, is a less strict metric that does not require the saliency method or benchmark annotators to locate the full extent of a pathology. Hit rate is based on the pointing game set-up⁴², in which credit is given if the most representative point identified by the saliency method or the benchmark annotators lies within the ground-truth segmentation. A ‘hit’ indicates that the correct region of the CXR was located regardless of the exact bounds of the binary segmentations. Localization performance is then calculated as the hit rate across the dataset⁴³. We report the means of these metrics (mIoU and hit rate) over 1,000 bootstrap replicates on the test set, along with the 95% confidence intervals using the 2.5th and 97.5th percentiles of the empirical distribution⁴⁴. In addition to mIoU, we report the test set precision, recall/sensitivity, and specificity values of the saliency method pipeline and the human benchmark segmentations to measure segmentation overlap (Extended Data Fig. 2).

Evaluating localization performance

To compare the localization performance of the saliency methods with the human benchmark, we first used Grad-CAM, Grad-CAM++ and integrated gradients to run 18 experiments, one for each combination of saliency method (Grad-CAM, Grad-CAM++ or integrated gradients) and CNN architecture (DenseNet121, ResNet152 or Inception-v4) using one of the two evaluation metrics (mIoU or hit rate) (Extended Data Fig. 3). We also ran experiments to evaluate the localization performances of DenseNet121 with Eigen-CAM, DeepLIFT, LRP and occlusion. We found that Grad-CAM with DenseNet121 generally demonstrated better localization performance across pathologies and evaluation metrics than the other combinations of saliency method and architecture. Accordingly, we compared Grad-CAM + DenseNet121 (saliency method pipeline) with the human benchmark using both mIoU and hit rate. See Table 1 for localization performance on the test set of all seven

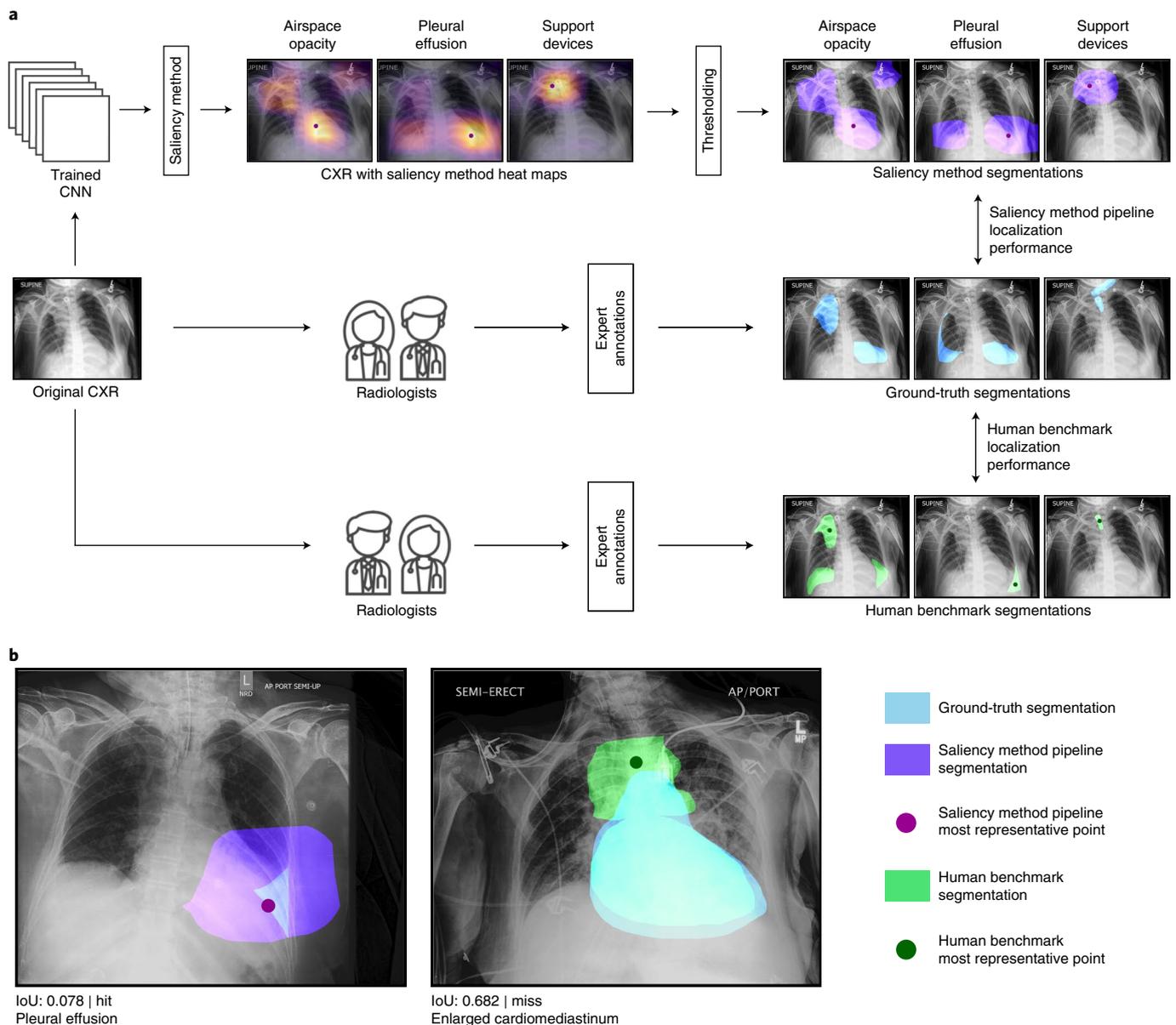


Fig. 1 | Framework for evaluating saliency methods. a, Top left: a CXR image from the holdout test set is passed into an ensemble CNN trained only on CXR images and their corresponding pathology task labels. The saliency method is used to generate ten heat maps for the example CXR, one for each task. The pixel in the heat map with the largest value is determined to be the single most representative point on the CXR for that pathology. There are three pathologies present in this CXR (airspace opacity, pleural effusion and support devices). Top right: a threshold is applied to the heat maps to produce binary segmentations for each present pathology. Middle row: two board-certified radiologists were asked to segment the pathologies that were present in the CXR as determined by the dataset's ground-truth labels. Saliency method pipeline annotations are compared with these ground-truth annotations to determine saliency method pipeline localization performance. Bottom row: three board-certified

radiologists (different from those of the middle row) were also asked to segment the pathologies that were present in the CXR as determined by the dataset's ground-truth labels. In addition, these radiologists were asked to locate the single point on the CXR that was most representative of each present pathology. These benchmark annotations are compared with the ground-truth annotations to determine human benchmark localization performance. **b**, Left: CXR with ground-truth and saliency method annotations for pleural effusion. The segmentations have a low overlap (IoU is 0.078), but the pointing game is a 'hit' since the saliency method's most representative point is inside the ground-truth segmentation. Right, CXR with ground-truth and human benchmark annotations for enlarged cardiomeastinum. The segmentations have a high overlap (IoU is 0.682), but the pointing game is a 'miss' since the saliency method's most representative point is outside the ground-truth segmentation.

saliency methods using DenseNet121. The localization performance for each pathology is reported on the true positive slice of the dataset (for mIoU, the true positive slice contains CXRs with both saliency method/human benchmark segmentations and also ground-truth segmentations; for hit rate, the true positive slice contains CXRs with both the most representative point identified by the saliency method/human benchmark and also the ground-truth segmentation). Localization

performance was calculated in this way so that saliency methods were not penalized by DNN classification error: while the benchmark radiologists were provided with ground-truth labels when annotating the dataset, saliency method segmentations were created on the basis of labels predicted by the model. (See Extended Data Fig. 4 for saliency method pipeline test set localization performance on the full dataset using mIoU.)

Table 1 | Test set localization performance of saliency methods using DenseNet121

Pathology	Grad-CAM	Grad-CAM++	Integrated gradients	Eigen-CAM	DeepLIFT	LRP	Occlusion
mIoU							
Airspace opacity	0.248	0.234	0.123	0.293	0.111	0.112	0.242
Atelectasis	0.254	0.245	0.116	0.267	0.126	0.109	0.250
Cardiomegaly	0.452	0.346	0.160	0.379	0.167	0.150	0.312
Consolidation	0.408	0.296	0.177	0.332	0.088	0.099	0.212
Edema	0.362	0.388	0.073	0.370	0.059	0.047	0.347
Enlarged cardiom.	0.379	0.400	0.154	0.372	0.109	0.117	0.363
Lung lesion	0.101	0.089	0.107	0.089	0.072	0.088	0.087
Pleural effusion	0.235	0.195	0.088	0.249	0.090	0.082	0.215
Pneumothorax	0.213	0.216	0.077	0.218	0.084	0.066	0.214
Support devices	0.163	0.133	0.099	0.116	0.086	0.052	0.126
Hit rate							
Airspace opacity	0.498	0.558	0.606	0.566	0.528	0.566	0.367
Atelectasis	0.501	0.621	0.520	0.530	0.415	0.468	0.343
Cardiomegaly	0.903	0.732	0.697	0.709	0.610	0.644	0.515
Consolidation	0.738	0.708	0.624	0.626	0.571	0.283	0.338
Edema	0.746	0.781	0.300	0.758	0.468	0.156	0.469
Enlarged cardiom.	0.818	0.630	0.704	0.612	0.469	0.594	0.767
Lung lesion	0.290	0.290	0.423	0.146	0.497	0.356	0.072
Pleural effusion	0.507	0.347	0.332	0.439	0.408	0.283	0.291
Pneumothorax	0.392	0.489	0.801	0.195	0.801	0.697	0.297
Support devices	0.355	0.364	0.491	0.216	0.598	0.264	0.189

We found that the saliency method pipeline demonstrated significantly worse localization performance on the test set when compared with the human benchmark using both mIoU (Fig. 2a) and hit rate (Fig. 2b) as an evaluation metric, regardless of the model classification AUROC (area under the receiver operating characteristic curve). For five of the ten pathologies, the saliency method pipeline had a significantly lower mIoU than the human benchmark. For example, the saliency method pipeline had one of the highest AUROC scores of the ten pathologies for support devices (0.969), but had among the worst localization performance for support devices when using both mIoU (0.163 [95% confidence interval (CI) 0.154, 0.172]) and hit rate (0.355 [95% CI 0.303, 0.408]) as evaluation metrics. On two pathologies (atelectasis and consolidation) the saliency method pipeline significantly outperformed the human benchmark. On average, across all ten pathologies, mIoU saliency method pipeline performance was 24.0% [95% CI 18.2%, 29.6%] worse than the human benchmark, with lung lesion displaying the largest gap in performance (76.2% [95% CI 59.1%, 87.5%] worse than the human benchmark) (Extended Data Fig. 5). Consolidation was the pathology on which the mIoU saliency method pipeline performance exceeded the human benchmark the most, by 128.1%. For seven of the ten pathologies, the saliency method pipeline had a significantly lower hit rate than the human benchmark. On average, hit rate saliency method pipeline performance was 29.4% [95% CI 23.1%, 35.5%] worse than the human benchmark (Extended Data Fig. 6), with lung lesion again displaying the largest gap in performance (65.9% [95% CI 35.3%, 91.7%] worse than the human benchmark). The hit rate saliency method pipeline did not significantly outperform the human benchmark on any of the ten pathologies; for the remaining three of the ten pathologies, the hit rate performance differences between the saliency method pipeline and the human benchmark were not statistically significant. Therefore, while the saliency method

pipeline significantly underperformed the human benchmark regardless of evaluation metric used, the average performance gap was larger when using hit rate as an evaluation metric than when using mIoU as an evaluation metric.

We compared saliency method pipeline localization performance using an ensemble model with localization performance using the top performing single checkpoint for each pathology. We found that on the test set the single model has worse localization performance than the ensemble model for all pathologies when using mIoU and for six of the ten pathologies when using hit rate (Extended Data Fig. 7).

Characterizing underperformance of saliency method pipeline

To better understand the underperformance of the saliency method pipeline localization, we first conducted a qualitative analysis with a radiologist by visually inspecting both the segmentations produced by the saliency method pipeline (Grad-CAM with DenseNet121) and the human benchmark segmentations. We found that, in general, saliency method segmentations fail to capture the geometric nuances of a given pathology, and instead produce coarse, low-resolution heat maps. Specifically, our qualitative analysis found that the performance of the saliency method was associated with three pathological characteristics (Fig. 3a): (1) number of instances (when a pathology had multiple instances on a CXR, the saliency method segmentation often highlighted one large confluent area, instead of highlighting each distinct instance of the pathology separately), (2) size (saliency method segmentations tended to be significantly larger than human expert segmentations, often failing to respect clear anatomical boundaries) and (3) shape complexity (the saliency method segmentations for pathologies with complex shapes frequently included significant portions of the CXR where the pathology is not present).

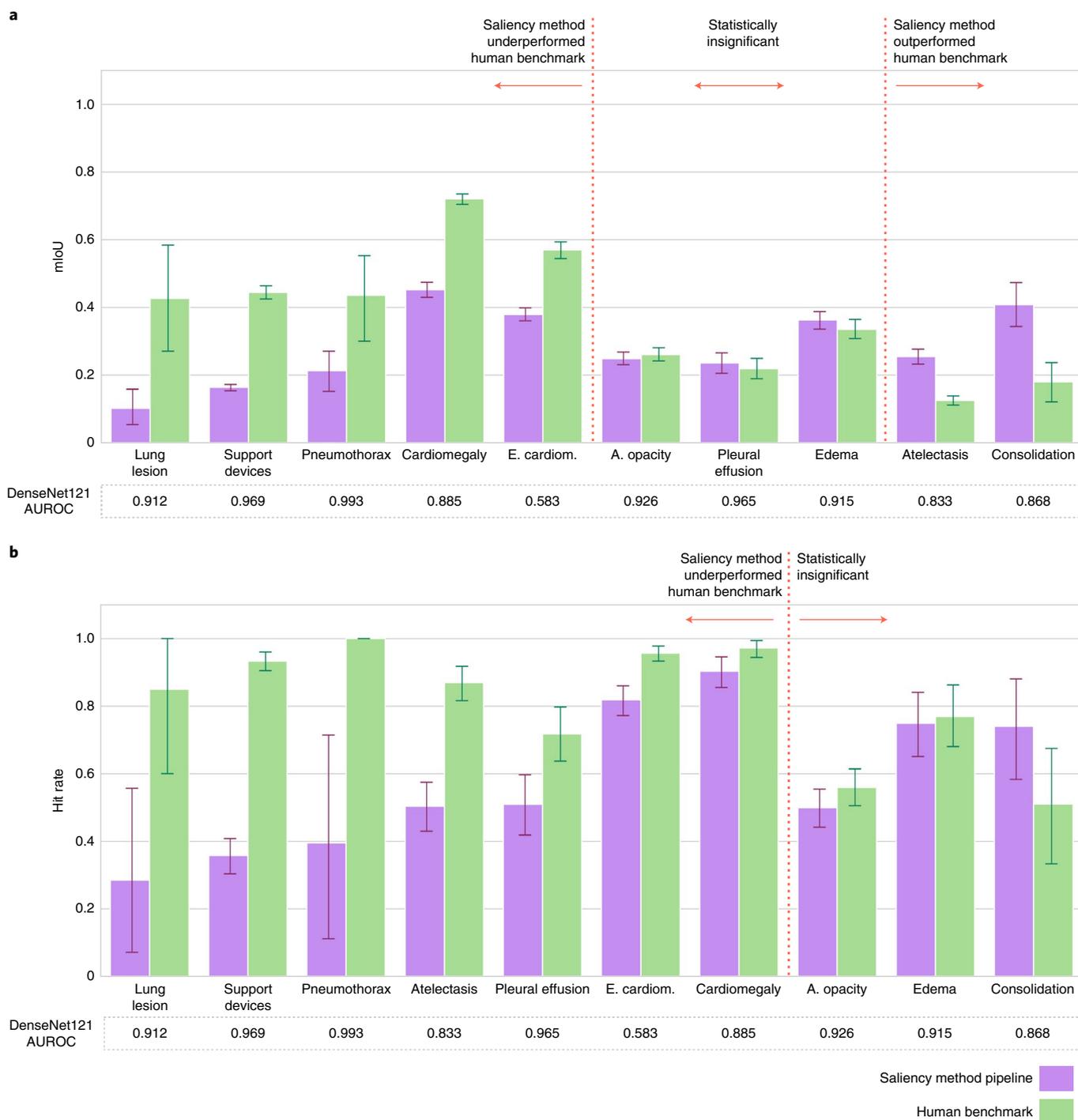


Fig. 2 | Evaluating localization performance. a, Comparing saliency method pipeline and human benchmark localization performances on the test set using mIoU. **b**, Comparing saliency method pipeline and human benchmark localization performances on the test set using hit rate. For both **a** and **b**, pathologies, along with their DenseNet121 AUROCs, are sorted on the x axis

first by statistical significance of percentage decrease from human benchmark mIoU/hit rate to saliency method pipeline mIoU/hit rate (high to low), and then by percentage decrease from human benchmark mIoU/hit rate to saliency method pipeline mIoU/hit rate (high to low).

Informed by our qualitative analysis and previous work in histology⁴⁵, we defined four geometric features for our quantitative analysis (Fig. 3b): (1) number of instances (for example, bilateral pleural effusion would have two instances, whereas there is only one instance for cardiomegaly), (2) size (pathology area with respect to the area of the whole CXR), (3) elongation and (4) irrectangularity (the last two features measure the complexity of the pathology shape and were

calculated by fitting a rectangle of minimum area enclosing the binary mask). See Extended Data Fig. 8 for the test set distribution of the four pathological characteristics across all ten pathologies.

For each evaluation metric, we ran eight simple linear regressions: four with the evaluation metric (IoU or hit/miss) of the saliency method pipeline (Grad-CAM with DenseNet121) as the dependent variable (to understand the relationship between the geometric features of a

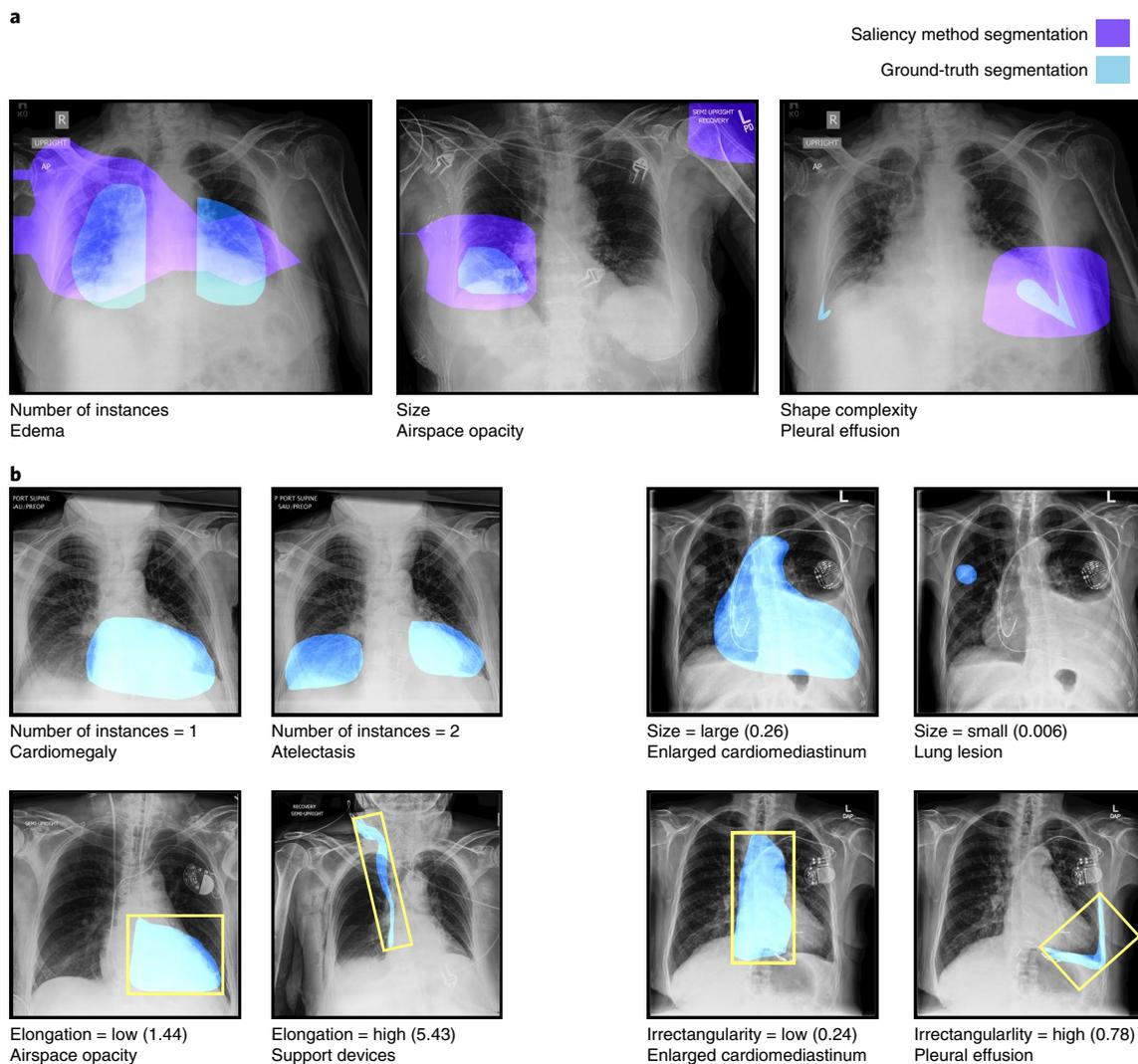


Fig. 3 | Characterizing underperformance of saliency method pipeline.

a, Example CXRs that highlight the three pathological characteristics identified by our qualitative analysis: (1) left, number of instances; (2) middle, size; (3) right, shape complexity. **b**, Example CXRs with the four geometric features used in our quantitative analysis: (1) top left, number of instances; (2) top right,

size = area of segmentation/area of CXR; (3) bottom left, elongation; (4) bottom right, irrectangularity. Elongation and irrectangularity were calculated by fitting a rectangle of minimum area enclosing the binary mask (as indicated by the yellow rectangles). Elongation = $\text{maxAxis}/\text{minAxis}$. Irrectangularity = $1 - (\text{area of segmentation}/\text{area of enclosing rectangle})$.

pathology and saliency method localization performance), and four with the difference between the evaluation metrics of the saliency method pipeline and the human benchmark as the dependent variable (to understand the relationship between the geometric features of a pathology and the gap in localization performance between the saliency method pipeline and the human benchmark). Each regression used one of the four geometric features as a single independent variable, and only the true positive slice was included in each regression. Each feature was normalized using min–max normalization and the regression coefficient can be interpreted as the effect of that geometric feature on the evaluation metric at hand. See Table 2 for coefficients from the regressions using both evaluation metrics on the test set, where we also report the 95% confidence interval and the Bonferroni-corrected *P* values based on Student's *t* distribution.

Our statistical analysis showed that as the size of a pathology increased the IoU saliency method localization performance improved (0.566 [95% CI 0.526, 0.606]). We also found that as elongation and irrectangularity increased the IoU saliency method localization performance worsened (elongation, -0.425 [95% CI -0.497 , -0.354]; irrectangularity, -0.256 [95% CI -0.292 , -0.219]). We observed that the

effects of these three geometric features were similar for hit/miss saliency method localization performance in terms of levels of statistical significance and direction of the effects. However, there was no evidence that the number of instances of a pathology had a significant effect on either IoU (-0.115 [95% CI -0.220 , -0.010]) or hit/miss (-0.051 [95% CI -0.346 , 0.244]) saliency method localization. Therefore, regardless of evaluation metric, saliency method localization performance suffered in the presence of pathologies that were small in size and complex in shape.

We found that these same three pathological characteristics—size, elongation and irrectangularity—characterized the gap in IoU localization performance between saliency method and human benchmark. We observed that the gap in hit/miss localization performance was significantly characterized by all four geometric features (number of instances, size, elongation and irrectangularity).

Effect of model confidence on localization performance

We also conducted statistical analyses to determine whether there was any correlation between the model's confidence in its prediction and saliency method pipeline test set localization performance (Table 3).

Table 2 | Coefficients from regressions on geometric features of pathologies

Geometric feature (independent variable)	Coefficient using saliency method localization (dependent variable)	Coefficient using localization difference (human benchmark – saliency method) (dependent variable)
IoU		
Number of instances	-0.115 (-0.220, -0.010)	-0.072 (-0.237, 0.094)
Size	0.566 (0.526, 0.606)***	-0.154 (-0.231, -0.076)***
Elongation	-0.425 (-0.497, -0.354)***	0.476 (0.362, 0.589)***
Irrectangularity	-0.256 (-0.292, -0.219)***	0.307 (0.249, 0.366)***
Hit/miss		
Number of instances	-0.051 (-0.346, 0.244)	0.470 (0.114, 0.825)*
Size	1.269 (1.146, 1.391)***	-0.944 (-1.104, -0.785)***
Elongation	-0.849 (-1.053, -0.646)***	1.110 (0.865, 1.354)***
Irrectangularity	-0.519 (-0.624, -0.415)***	0.689 (0.564, 0.815)***

* $P < 0.05$, *** $P < 0.001$.

We first ran a simple regression for each pathology using the model's probability output as the single independent variable and using the saliency method IoU as the dependent variable. We then performed a simple regression that uses the same approach as above, but includes all ten pathologies. For each of the 11 regressions, we used the full dataset since the analysis of false positives and false negatives was also of interest. In addition to the linear regression coefficients, we also computed the Spearman correlation coefficients to capture any potential nonlinear associations.

We found that for all pathologies the model confidence was positively correlated with the IoU saliency method pipeline performance. The P values for all coefficients were below 0.001 except for the coefficients for pneumothorax ($n = 11$) and lung lesion ($n = 50$), the two pathologies for which we had the fewest positive examples. Of all the pathologies, model confidence for positive predictions of enlarged cardiomeastinum had the largest linear regression coefficient with IoU saliency method pipeline performance (1.974, $P < 0.001$). Model confidence for positive predictions of pneumothorax had the largest Spearman correlation coefficient with IoU saliency method pipeline performance (0.734, $P < 0.05$), followed by pleural effusion (0.690, $P < 0.001$). Combining all pathologies ($n = 2,365$), the linear regression coefficient was 0.109 (95% CI [0.083, 0.135]), and the Spearman correlation coefficient was 0.285 (95% CI [0.248, 0.322]).

We also performed analogous experiments using hit/miss as the dependent variable on the true positive slice of the test set (CXRs with both the most representative point identified by the saliency method/human benchmark and also the ground-truth segmentations) (Extended Data Fig. 9). Since every heat map contains a maximally activated point (the pixel with the highest value) regardless of model probability output, using the full dataset has limited value since false positives are due to metric set-up and are not associated with model probability. We found that model confidence was positively correlated with hit/miss saliency method pipeline performance for four out of ten pathologies.

Discussion

The purpose of this work was to evaluate the performance of some of the most commonly used saliency methods for deep learning explainability using a variety of model architectures. We establish the first human benchmark for CXR segmentation in a multilabel classification set-up and demonstrate that saliency maps are consistently worse

Table 3 | IoU: coefficients from regressions on model assurance

Pathology	CXRs including all positives and false negatives (n)	Linear regression coefficient	Spearman correlation coefficient
Airspace opacity	381	0.714 (0.601, 0.826)***	0.577 (0.506, 0.641)***
Atelectasis	296	0.489 (0.333, 0.645)***	0.348 (0.244, 0.444)***
Cardiomegaly	229	0.679 (0.535, 0.823)***	0.592 (0.501, 0.670)***
Consolidation	120	1.155 (0.674, 1.635)***	0.384 (0.220, 0.527)***
Edema	124	0.642 (0.459, 0.826)***	0.548 (0.411, 0.660)***
Enlarged cardiomeastinum	668	1.974 (1.608, 2.340)***	0.428 (0.364, 0.488)***
Lung lesion	50	0.218 (0.088, 0.349)**	0.509 (0.268, 0.689)***
Pleural effusion	159	0.632 (0.489, 0.776)***	0.690 (0.599, 0.764)***
Pneumothorax	11	0.446 (0.108, 0.783)*	0.734 (0.240, 0.926)*
Support devices	327	0.211 (0.172, 0.250)***	0.468 (0.378, 0.548)***
All pathologies	2,365	0.109 (0.083, 0.135)***	0.285 (0.248, 0.322)***

* $P < 0.05$, ** $P < 0.01$, *** $P < 0.001$.

than expert radiologists regardless of model classification AUROC. We use qualitative and quantitative analyses to establish that saliency method localization performance is most inferior to expert localization performance when a pathology is smaller in size or has shapes that are more complex, suggesting that deep learning explainability as a clinical interface may be less reliable and less useful when used for pathologies with these characteristics. We also show that model assurance is positively correlated with saliency method localization performance, which could indicate that saliency methods are safer to use as a decision aid to clinicians when the model has made a positive prediction with high confidence.

Because ground-truth segmentations for medical imaging are time consuming and expensive to obtain, the current norm in medical imaging—both in research and in industry—is to use classification models on which saliency methods are applied post hoc for localization, highlighting the need for investigations into the reliability of these methods in clinical settings^{46,47}. There are public CXR datasets containing image-level labels annotated by expert radiologists (for example, the CheXpert validation set), multilabel bounding-box annotations (for example, ChestX-ray8⁴⁸ and VinDr-CXR⁴⁹) and segmentations for a single pathology (for example, SIIM-ACR pneumothorax segmentation⁵⁰). To our knowledge, however, there are no other publicly available CXR datasets with multilabel pixel-level expert segmentations. By publicly releasing a development dataset, CheXlocalize, of 234 images with 643 expert segmentations, we hope to encourage the further development of saliency methods and other explainability techniques for medical imaging.

Our work has several potential implications for human–AI (artificial intelligence) collaboration in the context of medical decision-making. Heat maps generated using saliency methods are advocated as clinical decision support in the hope that they not only improve clinical decision-making, but also encourage clinicians to trust model predictions^{51–53}. Many of the large CXR vendors (<https://annalise.ai/>, <https://www.lunit.io/en>, <https://qure.ai/>) use localization

methods to provide pathology visualization in their computer-aided detection products. In addition to being used for clinical interpretation, saliency method heat maps are also used for the evaluation of CXR interpretation models, for quality improvement and quality assurance in clinical practice, and for dataset annotation⁵⁴. Explainable AI is critical in high-stakes contexts such as healthcare, and saliency methods have been used successfully to develop and understand models generally. Indeed, we found that the saliency method pipeline significantly outperformed the human benchmark on two pathologies when using mIoU as an evaluation metric. However, our work also suggests that saliency methods are not yet reliable enough to validate individual clinical decisions made by a model. We found that saliency method localization performance, on balance, performed worse than expert localization across multiple analyses and across many important pathologies (our findings are consistent with recent work focused on localizing a single pathology, pneumothorax, in CXRs⁵⁵). We hypothesize that this could be an algorithmic artefact of saliency methods, whose relatively small heat maps (14 × 14 for Grad-CAM) are interpolated to the original image dimensions (usually 2,000 × 2,000), resulting in coarse resolutions. If used in clinical practice, heat maps that incorrectly highlight medical images may exacerbate well documented biases (chiefly automation bias) and erode trust in model predictions (even when model output is correct), limiting clinical translation²².

Since IoU computes the overlap of two segmentations but pointing game hit rate better captures diagnostic attention, we suggest using both metrics when evaluating localization performance in the context of medical imaging. While IoU is a commonly used metric for evaluating semantic segmentation outputs, there are inherent limitations to the metric in the pathological context. This is indicated by our finding that even the human benchmark segmentations had low overlap with the ground-truth segmentations (the highest expert mIoU was 0.720 for cardiomegaly). One potential explanation for this consistent underperformance is that pathologies can be hard to distinguish, especially without clinical context. Furthermore, whereas many people might agree on how to segment, say, a cat or a stop sign in traditional computer vision tasks, radiologists use a certain amount of clinical discretion when defining the boundaries of a pathology on a CXR. There can also be institutional and geographic differences in how radiologists are taught to recognize pathologies, and studies have shown that there can be high interobserver variability in the interpretation of CXRs^{56–58}. We sought to address this with the hit/miss evaluation metric, which highlights when two radiologists share the same diagnostic intention, even if it is less exact than IoU in comparing segmentations directly. The human benchmark localization using hit rate was above 0.9 for four pathologies (pneumothorax, cardiomegaly, enlarged cardiome-diastrinum and support devices); these are pathologies for which there is often little disagreement between radiologists about where the pathologies are located, even if the expert segmentations are noisy. Further work is needed to demonstrate which segmentation evaluation metrics, even beyond IoU and hit/miss, are more appropriate for certain pathologies and downstream tasks when evaluating saliency methods for the clinical setting.

Our work builds upon several studies investigating the validity of saliency maps for localization^{59–61} and upon some early work on the trustworthiness of saliency methods to explain DNNs in medical imaging⁴⁷. However, as recent work has shown³², evaluating saliency methods is inherently difficult given that they are post-hoc techniques. To illustrate this, consider the following models and saliency methods as described by some oracle: (1) a model *M_bad* that has perfect AUROC for a given image classification task, but that we know does not localize well (because the model picks up on confounders in the image); (2) a model *M_good* that also has perfect AUROC, but that we know does localize well (that is, is looking at relevant regions of the image); (3) a saliency method *S_bad* that does not properly reflect the model's attention; (4) a saliency method *S_good* that does properly reflect the model's

attention. Let us say that we are evaluating the following pipeline: we first classify an image and we then apply a saliency method post hoc. Imagine that our evaluation reveals poor localization performance as measured by mIoU or hit rate (as was the case in our findings). There are three possible pipelines (combinations of model and saliency method) that would lead to this scenario: (1) *M_bad* + *S_good*; (2) *M_good* + *S_bad*; (3) *M_bad* + *S_bad*. The first scenario (*M_bad* + *S_good*) is the one for which saliency methods were originally intended: we have a working saliency method that properly alerts us to models picking up on confounders. The second scenario (*M_good* + *S_bad*) is our nightmare scenario: we have a working model whose attention is appropriately directed, but we reject it on the basis of a poorly localizing saliency method. Because all three scenarios result in poor localization performance, it is difficult—if not impossible—to know whether poor localization performance is attributable to the model or to the saliency method (or to both). While we cannot say whether models or saliency methods are failing in the context of medical imaging, we can say that we should not rely on saliency methods to evaluate model localization. Future work should explore potential techniques for localization performance attribution.

There are several limitations of our work. First, we did not investigate the impact of pathology prevalence in the training data on saliency method localization performance. Second, some pathologies, such as effusions and cardiomegaly, are in similar locations across frontal view CXRs, while others, such as lesions and opacities, can vary in locations across CXRs. Future work could investigate how the locations of pathologies on a CXR in the training/test data distribution, and the consistency of these locations, affect saliency method localization performance. Third, while we compared saliency-method-generated pixel-level segmentations with human expert pixel-level segmentations, future work might explore how saliency method localization performance changes when comparing bounding-box annotations, instead of pixel-level segmentations. Fourth, we explored post-hoc interpretability methods given their prevalence in the context of medical imaging, but we hope that by publicly releasing our development dataset of pixel-level expert segmentations we can facilitate the development of models that make use of ground-truth segmentations during training⁵⁴. Fifth, the lack of a given finding can in certain cases inform clinical diagnoses. A common example of this is the lack of normal lung tissue pattern towards the edges of the thoracic cage, which is used to detect pneumothorax. For any characteristic pattern, both the absence and the presence provide diagnostic information to the radiologist. For example, the absence of a pleural effusion pattern is also used to rule out pleural effusion. For any characteristic radiological pattern, both the presence and the absence contribute to the final radiology report. Future work can explore counterfactual visual explanations that are similar to the counterfactual diagnostic process of a radiologist. Sixth, future work should further explore the potentially confounding effect of model calibration on the evaluation of saliency methods, especially when using segmentation, as opposed to classification, models. Finally, the impact of saliency methods on the trust and efficacy of users is underexplored.

In conclusion, we present a rigorous evaluation of a range of saliency methods and a dataset of pixel-level expert segmentations, which can serve as a foundation for future work exploring deep learning explainability techniques. This work is a reminder that care should be taken when leveraging common saliency methods to validate individual clinical decisions in deep learning-based workflows for medical imaging.

Methods

Ethical and information governance approvals

A formal review by the Stanford Institutional Review Board was conducted for the original collection of the CheXpert dataset. The Institutional Review Board waived the requirement to obtain

informed consent as the data were retrospectively collected and fully anonymized.

Dataset and clinical taxonomy. Dataset description. The localization experiments were performed using CheXpert, a large public dataset for CXR interpretation. The CheXpert dataset contains 224,316 CXRs for 65,240 patients labelled for the presence of 14 observations (13 pathologies and an observation of no finding) as positive, negative or uncertain. The CheXpert validation set consists of 234 CXRs from 200 patients randomly sampled from the full dataset and was labelled according to the consensus of three board-certified radiologists. The test set consists of 668 CXRs from 500 patients not included in the training or validation sets and was labelled according to the consensus of five board-certified radiologists. See Extended Data Fig. 10 for test set summary statistics. ‘Lung opacity’ in the CheXpert dataset is the equivalent of airspace opacity in the CheXlocalize dataset.

Ground-truth segmentations. The CXRs in our validation set and test set were manually segmented by two board-certified radiologists with 18 and 27 years of experience, using the annotation software tool MD.ai (<https://www.md.ai/>) (Supplementary Figs. 12–14). The radiologists were asked to contour the region of interest for all observations in the CXRs for which there was a positive ground-truth label in the CheXpert dataset. There were several cases in which the radiologists did not draw a certain pathology segmentation on a CXR even though the CXR had a positive ground-truth label for that pathology: airspace opacity on one CXR, atelectasis on one CXR, edema on two CXRs, enlarged cardiome-diastinum on one CXR and support devices on one CXR. For a pathology with multiple instances, all the instances were contoured. For support devices, radiologists were asked to contour any implanted or invasive devices (including pacemakers, peripherally inserted central catheters/central catheters, chest tubes, endotracheal tubes, feeding tubes and stents), and to ignore electrocardiography lead wires or external sticks visible in the CXR.

Benchmark segmentations. To evaluate expert performance on the test set using IoU, three radiologists, certified in Vietnam with 9, 10 and 18 years of experience, were asked to segment the regions of interest for all observations in the CXRs for which there was a positive ground-truth label in the CheXpert dataset. These radiologists were also provided with the same instructions for contouring as were provided to the radiologists drawing the ground-truth segmentations. To extract the maximally activated point from the benchmark segmentations, we asked the same radiologists to locate each pathology present on each CXR using only a single most representative point for that pathology on the CXR (see Supplementary Figs. 1–11 for the detailed instructions given to the radiologists). There was no overlap between these three radiologists and the two who drew the ground-truth segmentations.

Classification network architecture and training protocol. Multilabel classification model. The model takes as input a single-view CXR and outputs a probability for each of the 14 observations. If more than one view is available, the model outputs the maximum probability of the observations across the views. Each CXR was resized to 320×320 pixels and normalized before it was fed into the network. We used the same image resolutions as CheXpert⁴⁰ and CheXNet², which demonstrated radiologist-level performance on external test sets with 320×320 images. There are models that are commercially deployed and have similar dimensions. For example, the architecture used by medical AI software vendor Annalise.ai⁶² is based on EfficientNet⁶³, which takes input of 224×224 . CXRs were normalized before being fed into the network by subtracting the mean of all images in the CheXpert training set and then dividing by the s.d. of all images in the CheXpert training set. The model architectures DenseNet121, ResNet152 and Inception-v4 were used. Cross-entropy loss was used to train the model.

The Adam optimizer⁶⁴ was used with default β parameters of $\beta_1 = 0.9$ and $\beta_2 = 0.999$. The learning rate was tuned for the different model architectures using grid search (over learning rates of 1×10^{-3} , 1×10^{-4} and 1×10^{-5}). The best learning rate for each architecture was 1×10^{-4} for DenseNet121, 1×10^{-5} for ResNet152 and 1×10^{-5} for Inception-v4. Batches were sampled using a fixed batch size of 16 images.

Ensembling. We use an ensemble of checkpoints to create both predictions and saliency maps to maximize model performance. To capture uncertainties inherent in radiograph interpretation, we train our models using four uncertainty handling strategies outlined in CheXpert: ignoring, zeros, ones and three-class classification. For each of the four uncertainty handling strategies, we train our model three separate times, each time saving the ten checkpoints across the three epochs with the highest average AUROC across five observations selected for their clinical importance and prevalence in the validation set: atelectasis, cardiomegaly, consolidation, edema and pleural effusion. In total, after training, we have saved $4 \times 30 = 120$ checkpoints for a given model. Then, from the 120 saved checkpoints for that model, we select the ten top performing checkpoints for each pathology. For each CXR and each task, we compute the predictions and saliency maps using the relevant checkpoints. We then take the mean both of the predictions and of the saliency maps to create the final set of predictions and saliency maps for the ensemble model. See Supplementary Table 1 for the performance of each model architecture (DenseNet121, ResNet152 and Inception-v4) with each of the pathologies.

Evaluating localization performance. Saliency methods were used to visualize the decision made by the classification network. Each saliency map was resized to the original image dimension using bilinear interpolation. It was then normalized using min–max normalization and converted into a binary segmentation using binary thresholding (Otsu’s method). For occlusion, we used a window size of 40 and a stride of 40 for each CXR.

Localization performance of each segmentation was evaluated using IoU score. The IoU is the ratio between the area of overlap and the area of union between the ground-truth and the predicted segmentations, ranging from 0 to 1 (0 signifies no overlap and 1 signifies perfectly overlapping segmentations). We report the mIoU over 1,000 bootstrap replicates on the test set, along with the 95% confidence intervals using the 2.5th and 97.5th percentiles of the empirical distribution.

For the evaluation of DenseNet121 + integrated gradients using IoU, we applied box filtering of kernel size 100 to smooth the pixelated map. For the evaluation of ResNet152 + integrated gradients and of Inception-v4 + integrated gradients using IoU, we applied box filtering of kernel size 50. For the evaluation of DeepLIFT using IoU, we applied box filtering of kernel size 50. For the evaluation of LRP using IoU, we applied box filtering of kernel size 80. The kernel sizes were tuned on the validation set. The noisy map is not a concern for hit rate because a single maximum pixel is extracted for the entire image.

In Extended Data Fig. 1, we report mIoU localization performance using different saliency map thresholding values. We first applied min–max normalizations to the saliency maps so that each value is transformed into a decimal between 0 and 1. We then passed in a range of threshold values from 0.2 to 0.8 to create binary segmentations and calculated the mIoU score per pathology under each threshold on the validation set.

In Extended Data Fig. 2, we report the precision, recall/sensitivity, and specificity values of the saliency method pipeline and the human benchmark segmentations on the test set.

For this, we treat each pixel in the saliency method pipeline and the human benchmark segmentations as a classification, use each pixel in the ground-truth segmentation as the ground-truth label, and calculate precision, recall/sensitivity, and specificity over all CXRs for each pathology. Precision is defined as total number of true positive

pixels/(total number of true positive + false positive pixels). Recall is defined as total number of true positive pixels/(total number of true positive + false negative pixels). Specificity is defined as total number of true negative pixels/(total number of true negative + false positive pixels).

In Extended Data Fig. 4, we report the saliency method pipeline test set localization performance on the full dataset using mIoU. For this, we ensure that the final binary segmentation is consistent with model probability output by applying another layer of thresholding such that the segmentation mask produces all zeros if the predicted probability is below a chosen level. The probability threshold is searched on the interval of [0, 0.8] with steps of 0.1. The exact value is determined per pathology by maximizing the mIoU on the validation set.

In Extended Data Figs. 5 and 6, we report the percentage decrease from human benchmark localization performance to saliency method pipeline localization performance on the test set. To obtain the 95% confidence interval per pathology on the percentage decrease from human benchmark localization performance to saliency method pipeline localization performance, we first extracted the percentage decrease statistic $([\text{human benchmark mIoU or hit rate} - \text{saliency method pipeline mIoU or hit rate}] / \text{human benchmark mIoU or hit rate} \times 100)$ from each of the 1,000 human benchmark and the 1,000 saliency method pipeline mIoU/hit rate bootstrap replicates for each pathology. In doing so, we created the bootstrap distribution of the percentage decrease statistic. We reported the 95% CI using the 2.5th and 97.5th percentiles of the empirical distribution. To obtain the 95% CI on the average percentage decrease over all pathologies, the methodology is the same: we created bootstrap replicates of the average human benchmark and saliency method pipeline mIoUs/hit rates over all pathologies, extracted the percentage decrease statistic from each replicate and then reported the 95% CI using the 2.5th and 97.5th percentiles of the empirical distribution.

Statistical analyses

Pathology characteristics. The pathology characteristics used in all regressions were calculated on the ground-truth annotations. The four characteristics are defined as follows. (1) Number of instances is the number of separate segmentations drawn by the radiologist for a given pathology. (2) Size is the area of the pathology divided by the total image area. (3), (4) Elongation and irrectangularity are geometric features that measure shape complexities. They were designed to quantify what radiologists qualitatively describe as focal or diffused. To calculate the metrics, a rectangle of minimum area enclosing the contour is fitted to each pathology. Elongation is defined as the ratio of the rectangle's longer side to shorter side. Irrectangularity = $1 - (\text{area of segmentation/area of enclosing rectangle})$, with values ranging from 0 to 1 (1 being very irrectangular). When there were multiple instances within one pathology, we used the characteristics of the dominant instance (largest in perimeter). All geometric features are normalized using min–max normalization per pathology before aggregation so that they are comparable on scales of magnitudes.

Model confidence. We used the probability output of the DNN architecture for model confidence. The probabilities were on a similar scale of 0–1 and we did not apply min–max normalization. We report the 95% confidence interval and *P* value of the regression coefficients using Student's *t* distribution.

For the statistical analyses on the full dataset to determine whether there was any correlation between the model's confidence in its prediction and saliency method pipeline performance using IoU (Table 3), we ensured that the final binary segmentation was consistent with model probability output by applying another layer of thresholding such that the segmentation mask produced all zeros if the predicted probability was below a chosen level. The probability threshold is searched on the interval of [0, 0.8] with steps of 0.1. The exact value is determined per pathology by maximizing the mIoU on the validation set.

For the statistical analyses to determine whether there was any correlation between the model's confidence in its prediction and saliency method pipeline performance using hit/miss (Extended Data Fig. 9), we used the true positive slice of the dataset (CXRs with both the most representative point identified by the saliency method/human benchmark and also the ground-truth segmentation). Since every heat map contains a maximally activated point (the pixel with the highest value) regardless of model probability output, using the full dataset has limited value since false positives are due to metric set-up and are not associated with model probability.

Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

The CheXlocalize dataset is available here: <https://stanfordaimi.azurewebsites.net/datasets/abfb76e5-70d5-4315-badc-c94dd82e3d6d>. The CheXpert dataset is available here <https://stanfordmlgroup.github.io/competitions/chexpert/>.

Code availability

The code used to produce our results is available in the following public repository under the MIT License: <https://github.com/rajpurkar-lab/cheXlocalize>. The version used for this publication is available at <https://doi.org/10.5281/zenodo.6973536>.

References

- Rajpurkar, P. et al. Deep learning for chest radiograph diagnosis: a retrospective comparison of the CheXNeXt algorithm to practicing radiologists. *PLoS Med.* **15**, e1002686 (2018).
- Rajpurkar, P. et al. CheXNet: radiologist-level pneumonia detection on chest X-rays with deep learning. Preprint at: <https://arxiv.org/abs/1711.05225> (2017).
- Bien, N. et al. Deep-learning-assisted diagnosis for knee magnetic resonance imaging: development and retrospective validation of MRNet. *PLoS Med.* **15**, e1002699 (2018).
- Baselli, G., Codari, M. & Sardanelli, F. Opening the *black box* of machine learning in radiology: can the proximity of annotated cases be a way? *Eur. Radiol. Exp.* **4**, 30 (2020).
- Litjens, G. et al. A survey on deep learning in medical image analysis. *Med. Image Anal.* **42**, 60–88 (2017).
- Wang, F., Kaushal, R. & Khullar, D. Should health care demand interpretable artificial intelligence or accept 'black box' medicine? *Ann. Intern. Med.* **172**, 59–60 (2019).
- Goodman, B. & Flaxman, S. European Union regulations on algorithmic decision-making and a 'right to explanation'. *AI Mag.* **38**, 50–57 (2017).
- Venugopal, V. K., Takhar, R., Gupta, S., Saboo, A. & Mahajan, V. Clinical Explainability Failure (CEF) & Explainability Failure Ratio (EFR)—changing the way we validate classification algorithms? *J. Med. Syst.* **46**, 20 (2022).
- Pasa, F., Golkov, V., Pfeiffer, F., Cremers, D. & Pfeiffer, D. Efficient deep network architectures for fast chest X-ray tuberculosis screening and visualization. *Sci. Rep.* **9**, 6268 (2019).
- Simonyan, K., Vedaldi, A. & Zisserman, A. Deep inside convolutional networks: visualising image classification models and saliency maps. Workshop at International Conference on Learning Representations (2014).
- Aggarwal, M. et al. Towards trainable saliency maps in medical imaging. Machine Learning for Health (ML4H) Extended Abstract Arxiv, Index:1–6 (2020).
- Tjoa, E. & Guan, C. Quantifying explainability of saliency methods in deep neural networks. Preprint at: <https://arxiv.org/abs/2009.02899> (2020).

13. Badgeley, M. A. et al. Deep learning predicts hip fracture using confounding patient and healthcare variables. *npj Digit. Med.* **2**, 31 (2019).
14. Zech, J. R. et al. Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: a cross-sectional study. *PLoS Med.* **15**, e1002683 (2018).
15. DeGrave, A. J., Janizek, J. D. & Lee, S.-I. AI for radiographic COVID-19 detection selects shortcuts over signal. *Nat. Mach. Intell.* **3**, 610–619 (2021).
16. Makimoto, H. et al. Performance of a convolutional neural network derived from an ECG database in recognizing myocardial infarction. *Sci. Rep.* **10**, 8445 (2020).
17. Porumb, M., Stranges, S., Pescapè, A. & Pecchia, L. Precision medicine and artificial intelligence: a pilot study on deep learning for hypoglycemic events detection based on ECG. *Sci. Rep.* **10**, 170 (2020).
18. Tham, Y.-C. et al. Referral for disease-related visual impairment using retinal photograph-based deep learning: a proof-of-concept, model development study. *Lancet Digit. Health* **3**, e29–e40 (2021).
19. Varadarajan, A. V. et al. Deep learning for predicting refractive error from retinal fundus images. *Invest. Ophthalmol. Vis. Sci.* **59**, 2861–2868 (2018).
20. Mitani, A. et al. Detection of anaemia from retinal fundus images via deep learning. *Nat. Biomed. Eng.* **4**, 18–27 (2020).
21. Lu, M. T. et al. Deep learning to assess long-term mortality from chest radiographs. *JAMA Netw. Open* **2**, e197416 (2019).
22. Rajpurkar, P. et al. CheXaid: deep learning assistance for physician diagnosis of tuberculosis using chest x-rays in patients with HIV. *npj Digit. Med.* **3**, 115 (2020).
23. Rajpurkar, P. et al. AppendixNet: deep learning for diagnosis of appendicitis from a small dataset of CT exams using video pretraining. *Sci. Rep.* **10**, 3958 (2020).
24. Huang, S.-C. et al. PENet—a scalable deep-learning model for automated diagnosis of pulmonary embolism using volumetric CT imaging. *npj Digit. Med.* **3**, 61 (2020).
25. Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat. Mach. Intell.* **1**, 206–215 (2019).
26. Eitel, F. et al. Testing the robustness of attribution methods for convolutional neural networks in MRI-based Alzheimer’s disease classification. In *Interpretability of Machine Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision Support. ML-CDS IMIMIC 2019* (eds Suzuki, K. et al.) 3–11 (Lecture Notes in Computer Science Vol. 11797, Springer, 2019).
27. Young, K., et al. Deep neural network or dermatologist? In *Interpretability of Machine Intelligence in Medical Image Computing and Multimodal Learning for Clinical Decision Support. ML-CDS IMIMIC 2019* (eds Suzuki, K. et al.) 48–55 (Lecture Notes in Computer Science Vol. 11797, Springer, 2019).
28. Ghassemi, M., Oakden-Rayner, L. & Beam, A. L. The false hope of current approaches to explainable artificial intelligence in health care. *Lancet Digit. Health* **3**, e745–e750 (2021).
29. Reyes, M. et al. On the interpretability of artificial intelligence in radiology: challenges and opportunities. *Radiol. Artif. Intell.* **2**, e190043 (2020).
30. Selvaraju, R. R. et al. Grad-CAM: visual explanations from deep networks via gradient-based localization. *Int. J. Comput. Vis.* **128**, 336–359 (2020).
31. Chattopadhyay, A., Sarkar, A., Howlader, P. & Balasubramanian, V. N. Grad-CAM++: generalized gradient-based visual explanations for deep convolutional networks. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)* 839–847 (IEEE, 2018).
32. Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. *Proc. Mach. Learning Res.* **70**, 3319–3328 (2017).
33. Bany Muhammad, M. et al. Eigen-CAM: visual explanations for deep convolutional neural networks. *SN Comput. Sci.* **2**, 47 (2021).
34. Shrikumar, A., Greenside, P. & Kundaje, A. Learning important features through propagating activation differences. *Proc. Mach. Learning Res.* **70**, 3145–3153 (2017).
35. Bach, S. et al. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE* **10**, e0130140 (2015).
36. Zeiler, M. D. & Fergus, R. Visualizing and understanding convolutional networks. In *Computer Vision—ECCV 2014* (eds Fleet, D. et al.) 818–833 (Springer, 2014).
37. Huang, G., Liu, Z., Van Der Maaten, L. & Weinberger, K. Q. Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 2261–2269 (IEEE, 2017).
38. He, K., Zhang, X., Ren, S. & Sun, J. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 770–778 (IEEE, 2016).
39. Szegedy, C. et al. Going deeper with convolutions. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 1–9 (IEEE, 2015).
40. Irvin, J. et al. CheXpert: a large chest radiograph dataset with uncertainty labels and expert comparison. In *Proc. AAAI Conference on Artificial Intelligence* Vol. 33, 590–597 (AAAI, 2019).
41. Otsu, N. A threshold selection method from gray-level histograms. *IEEE Trans. Syst. Man Cybern.* **9**, 62–66 (1979).
42. Zhang, J. et al. Top-down neural attention by excitation backprop. *Int. J. Comput. Vis.* **126**, 1084–1102 (2018).
43. Kim, H.-E. et al. Changes in cancer detection and false-positive recall in mammography using artificial intelligence: a retrospective, multireader study. *Lancet Digit. Health* **2**, e138–e148 (2020).
44. Efron, B. & Tibshirani, R. J. *An Introduction to the Bootstrap* (CRC Press, 1994).
45. Vrabac, D. et al. DLBCL-Morph: morphological features computed using deep learning for an annotated digital DLBCL image set. *Sci. Data* **8**, 135 (2021).
46. Ayhan, M. S. et al. Clinical validation of saliency maps for understanding deep neural networks in ophthalmology. *Med. Image Anal.* **77**, 102364 (2022).
47. Arun, N. et al. Assessing the trustworthiness of saliency maps for localizing abnormalities in medical imaging. *Radiol. Artif. Intell.* **3**, e200267 (2021).
48. Wang, X. et al. ChestX-ray8: hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 2097–2106 (IEEE, 2017).
49. Nguyen, H. Q. et al. VinDr-CXR: An open dataset of chest X-rays with radiologist’s annotations. *Sci. Data* **9**, 429 (2022).
50. Society for Imaging Informatics in Medicine (SIIM) SIIM-ACR pneumothorax segmentation. Kaggle <https://kaggle.com/c/siim-acr-pneumothorax-segmentation> (2019).
51. Steiner, D. F. et al. Impact of deep learning assistance on the histopathologic review of lymph nodes for metastatic breast cancer. *Am. J. Surg. Pathol.* **42**, 1636–1646 (2018).
52. Uyumazturk, B. et al. Deep learning for the digital pathologic diagnosis of cholangiocarcinoma and hepatocellular carcinoma: evaluating the impact of a web-based diagnostic assistant. Machine Learning for Health (ML4H) at NeurIPS - Extended Abstract (2019).
53. Park, A. et al. Deep learning-assisted diagnosis of cerebral aneurysms using the HeadXNet model. *JAMA Netw. Open* **2**, e195600 (2019).
54. Gadgil, S., Endo, M., Wen, E., Ng, A. Y. & Rajpurkar, P. CheXseg: combining expert annotations with DNN-generated saliency

- maps for X-ray segmentation. *Proc. Mach. Learning Res.* **143**, 190–204 (2021).
55. Crosby, J., Chen, S., Li, F., MacMahon, H. & Giger, M. Network output visualization to uncover limitations of deep learning detection of pneumothorax. *Proc. SPIE* **11316**, 113160O (2020).
 56. Melbye, H. & Dale, K. Interobserver variability in the radiographic diagnosis of adult outpatient pneumonia. *Acta Radiol.* **33**, 79–81 (1992).
 57. Herman, P. G. et al. Disagreements in chest Roentgen interpretation. *CHEST* **68**, 278–282 (1975).
 58. Albaum, M. N. et al. Interobserver reliability of the chest radiograph in community-acquired pneumonia. *CHEST* **110**, 343–350 (1996).
 59. Arun, N. T. et al. Assessing the validity of saliency maps for abnormality localization in medical imaging. In Tal Arbel, Ben Ayed, Marleen de Bruijne, Maxime Descoteaux, Herve Lombaert, Chris Pal (eds.), *Medical Imaging with Deep Learning 2020*, Short Paper Track (2020).
 60. Graziani, M., Lompech, T., Müller, H. & Andrearczyk, V. Evaluation and comparison of CNN visual explanations for histopathology. In *AAAI 2021, XAI Workshop* (2021).
 61. Choe, J. et al. Evaluating weakly supervised object localization methods right. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition* 3133–3142 (IEEE, 2020).
 62. Seah, J. C. Y. et al. Effect of a comprehensive deep-learning model on the accuracy of chest x-ray interpretation by radiologists: a retrospective, multireader multicase study. *Lancet Digit. Health* **3**, e496–e506 (2021).
 63. Tan, M. & Le, Q. V. EfficientNet: rethinking model scaling for convolutional neural networks. *Proc. Mach. Learning Res.* **97**, 6105–6114 (2019).
 64. Kingma, D. P. & Ba, J. Adam: a method for stochastic optimization. International Conference on Learning Representations (ICLR) Poster (2015).
 65. Saporta, A. et al. Code for ‘Benchmarking saliency methods for chest X-ray interpretation’. *Zenodo* <https://doi.org/10.5281/zenodo.6973536> (2022).

Acknowledgements

We acknowledge MD.ai for providing us access to their annotation platform. We acknowledge Weights & Biases for providing us access to their experiment tracking tools.

Author contributions

Conceptualization: P.R. and A.P. Design: P.R., A.P., A.S., X.G. and A.A. Data analysis and interpretation: A.S., X.G., A.A., P.R., A.P., S.Q.H.T., C.D.T.N., V.-D.N., J.S. and F.G.B. Drafting of the manuscript: A.S., X.G., A.A. and P.R. Critical revision of the manuscript for important

intellectual content: A.P., S.Q.H.T., C.D.T.N., V.-D.N., J.S., F.G.B., A.Y.N. and M.P.L. Supervision: A.Y.N., M.P.L. and P.R. Research was primarily performed while A.S. was at Stanford University. M.P.L. and P.R. contributed equally.

Competing interests

M.P.L. is an adviser for and/or has research funded by GE, Philips, Carestream, Nines Radiology, Segmed, Centaur Labs, Microsoft, Bunkerhill and Amazon Web Services (none of the funded research was relevant to this project). A.P. is a medical associate at Cerebriu. The remaining authors declare no competing interests.

Additional information

Extended data is available for this paper at <https://doi.org/10.1038/s42256-022-00536-x>.

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s42256-022-00536-x>.

Correspondence and requests for materials should be addressed to Pranav Rajpurkar.

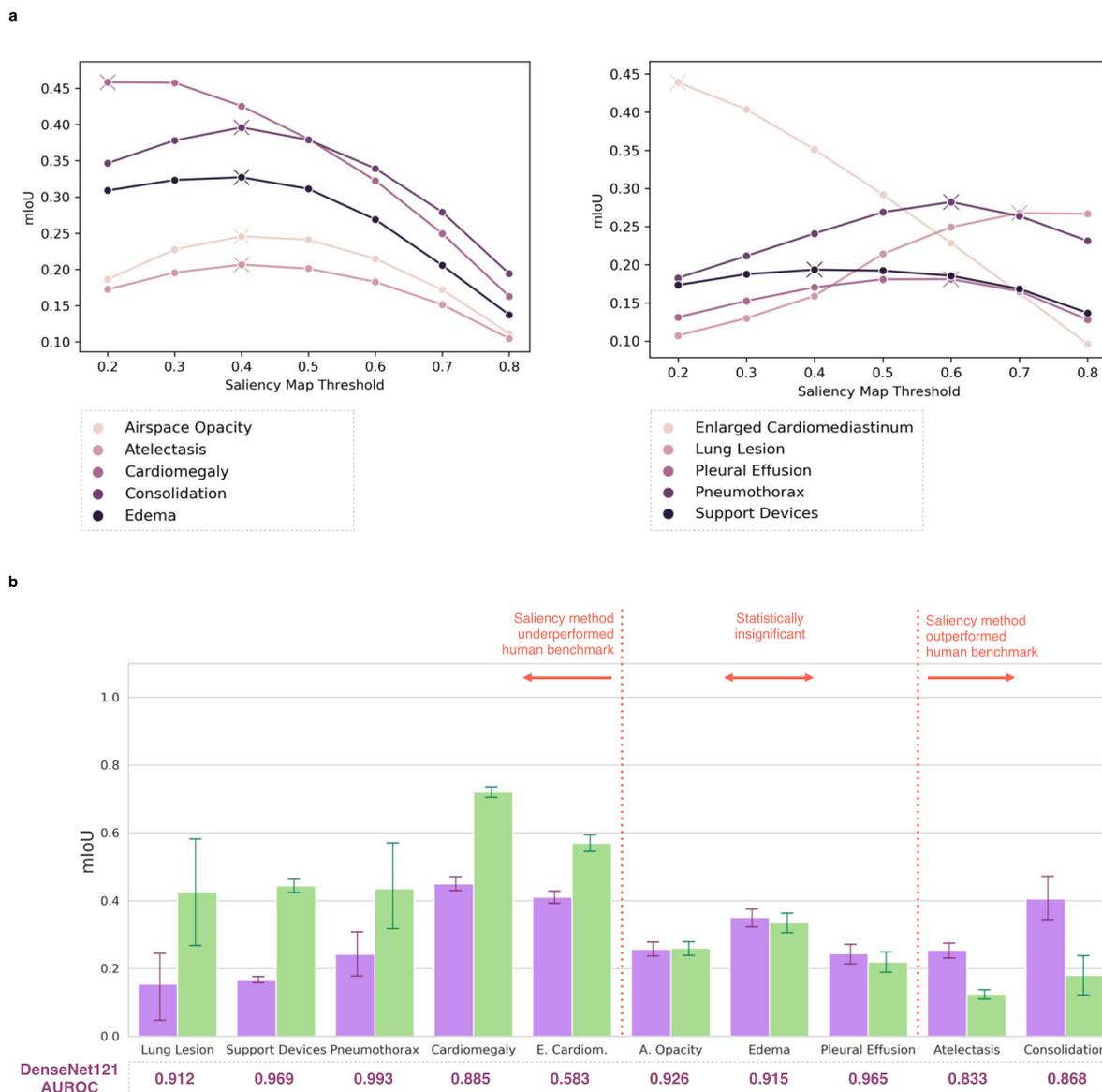
Peer review information *Nature Machine Intelligence* thanks Alex Zwanenburg, Xiaosong Wang and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2022



Extended Data Fig. 1 | mIoU localization performance of the saliency method pipeline on the test set using threshold values tuned on the validation set.

a, We first applied min-max normalization to the Grad-CAM saliency maps so that each value gets transformed into a decimal between 0 and 1. We then passed in a range of threshold values from 0.2 to 0.8 to create binary segmentations and plotted the mIoU score per pathology under each threshold on the validation set. The threshold that gives the max mIoU for each pathology is marked with an “X”. Pathologies are sorted alphabetically and shown in two plots for readability.

b, Comparing mIoU localization performances of the saliency method pipeline on the test set (using the best thresholds tuned on the validation set) and the human benchmark. We found that the saliency method pipeline outperformed the human benchmark on two pathologies and underperformed the human benchmark on five pathologies. For the remaining three pathologies, the performance differences were not statistically significant. This finding is consistent with what we report in the manuscript using Otsu’s method.

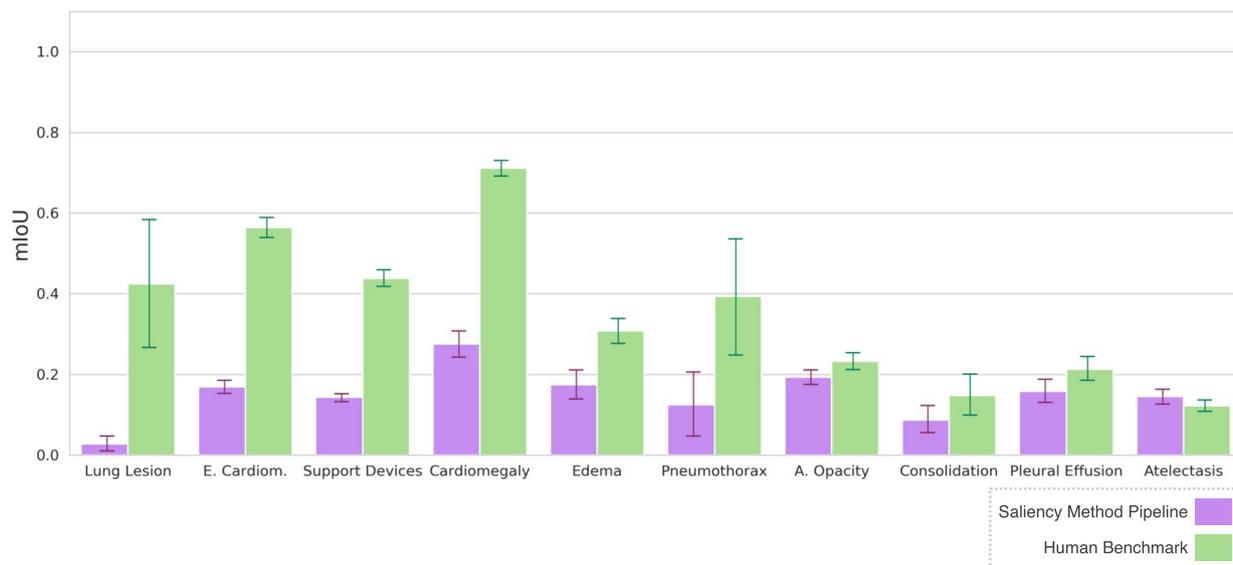
pathology	precision		recall/sensitivity		specificity	
	saliency method pipeline	human benchmark	saliency method pipeline	human benchmark	saliency method pipeline	human benchmark
Airspace Opacity	0.194	0.526	0.638	0.328	0.838	0.982
Atelectasis	0.119	0.838	0.593	0.111	0.852	0.999
Cardiomegaly	0.277	0.947	0.621	0.728	0.911	0.998
Consolidation	0.045	0.701	0.639	0.158	0.915	1.000
Edema	0.114	0.753	0.652	0.330	0.848	0.997
Enlarged Cardiom.	0.454	0.905	0.467	0.622	0.931	0.992
Lung Lesion	0.003	0.614	0.797	0.469	0.887	1.000
Pleural Effusion	0.052	0.653	0.673	0.154	0.806	0.999
Pneumothorax	0.005	0.838	0.786	0.459	0.866	1.000
Support Devices	0.120	0.849	0.510	0.425	0.860	0.997

Extended Data Fig. 2 | Precision, recall/sensitivity, and specificity values of the saliency method pipeline and the human benchmark segmentations on the test set. We treated each pixel in the saliency method pipeline and the human benchmark segmentations as a classification, used each pixel in the ground-truth

segmentation as the ground-truth label, and calculated the precision, recall/sensitivity, and specificity over all CXRs for each pathology. For each pathology and each metric, we highlight the higher of the two (saliency method pipeline or human benchmark) in **bold**.

pathology	Grad-CAM			Grad-CAM++			Integrated Gradients		
	DenseNet121	ResNet152	Inception-v4	DenseNet121	ResNet152	Inception-v4	DenseNet121	ResNet152	Inception-v4
mIoU									
Airspace Opacity	0.248	0.194	0.090	0.234	0.198	0.115	0.123	0.119	0.052
Atelectasis	0.254	0.221	0.115	0.245	0.210	0.106	0.116	0.115	0.064
Cardiomegaly	0.452	0.424	0.120	0.346	0.257	0.196	0.160	0.154	0.089
Consolidation	0.408	0.334	0.079	0.296	0.245	0.130	0.177	0.112	0.069
Edema	0.362	0.240	0.203	0.388	0.345	0.266	0.073	0.062	0.099
Enlarged Cardiom.	0.379	0.272	0.065	0.400	0.382	0.295	0.154	0.152	0.094
Lung Lesion	0.101	0.066	0.003	0.089	0.069	0.045	0.107	0.063	0.001
Pleural Effusion	0.235	0.204	0.120	0.195	0.176	0.090	0.088	0.091	0.067
Pneumothorax	0.213	0.171	0.088	0.216	0.184	0.124	0.077	0.070	0.078
Support Devices	0.163	0.147	0.116	0.133	0.126	0.099	0.099	0.074	0.066
hit rate									
Airspace Opacity	0.498	0.428	0.106	0.558	0.522	0.148	0.606	0.586	0.122
Atelectasis	0.501	0.490	0.062	0.621	0.621	0.118	0.520	0.453	0.187
Cardiomegaly	0.903	0.915	0.126	0.732	0.297	0.493	0.697	0.748	0.268
Consolidation	0.738	0.797	0.030	0.708	0.600	0.284	0.624	0.538	0.115
Edema	0.746	0.432	0.385	0.781	0.745	0.457	0.300	0.350	0.180
Enlarged Cardiom.	0.818	0.627	0.030	0.630	0.631	0.731	0.704	0.730	0.205
Lung Lesion	0.290	0.146	0.000	0.290	0.146	0.000	0.423	0.211	0.000
Pleural Effusion	0.507	0.499	0.133	0.347	0.473	0.107	0.332	0.400	0.182
Pneumothorax	0.392	0.600	0.000	0.489	0.698	0.097	0.801	0.498	0.097
Support Devices	0.355	0.287	0.133	0.364	0.334	0.150	0.491	0.442	0.324

Extended Data Fig. 3 | Test set localization performance for each combination of saliency method and CNN architecture. For each pathology and saliency method, we highlight the highest performing CNN architecture in bold.



Extended Data Fig. 4 | Saliency method pipeline test set localization performance on the full dataset using mIoU. True negatives (CXRs whose ground-truth label is negative for a given pathology and for which there were neither human benchmark nor saliency method pipeline segmentations for that pathology) were excluded from the metric calculation. To control for false positives, we ensure that the final binary segmentation is consistent with model probability output by applying another layer of thresholding such that the

segmentation mask produced all zeros if the predicted probability was below a chosen level. The probability threshold is searched on the interval of $[0, 0.8]$ with steps of 0.1. The exact value is determined per pathology by maximizing the mIoU on the validation set. We found that on the full dataset, for seven of the 10 pathologies, the saliency method pipeline had a significantly lower mIoU than the human benchmark.

pathology	human benchmark mIoU	saliency method pipeline mIoU	% decrease (95% CI)
Lung Lesion	0.426	0.101	76.2 (59.1, 87.5)
Support Devices	0.444	0.163	63.3 (60.8, 65.8)
Pneumothorax	0.435	0.213	50.9 (14.6, 69.5)
Cardiomegaly	0.720	0.452	37.2 (34.0, 40.4)
Enlarged Cardiom.	0.569	0.379	33.4 (29.0, 37.4)
Airspace Opacity	0.260	0.248	4.8 (-6.1, 14.6)
Pleural Effusion	0.219	0.235	-7.6 (-34.5, 13.3)
Edema	0.335	0.362	-7.9 (-19.7, 2.6)
Atelectasis	0.124	0.254	-104.4 (-134.2, -78.2)
Consolidation	0.179	0.408	-128.1 (-226.8, -74.5)
Average	0.371	0.282	24.0 (18.2, 29.6)

Extended Data Fig. 5 | Percentage decrease from human benchmark mIoU to saliency method pipeline mIoU on the test set. Pathologies are sorted first by statistical significance of percentage decrease from human benchmark mIoU to

saliency method pipeline mIoU (high to low), and then by percentage decrease from human benchmark mIoU to saliency method pipeline mIoU (high to low). We use 95% bootstrap confidence interval.

pathology	human benchmark hit rate (%)	saliency method pipeline hit rate (%)	% decrease (95% CI)
Lung Lesion	0.850	0.290	65.9 (35.3, 91.7)
Support Devices	0.933	0.355	61.9 (56.2, 67.5)
Pneumothorax	1.000	0.392	60.8 (27.3, 92.3)
Atelectasis	0.870	0.501	42.4 (33.1, 51.0)
Pleural Effusion	0.718	0.507	29.4 (14.3, 42.5)
Enlarged Cardiom.	0.957	0.818	14.5 (9.6, 19.2)
Cardiomegaly	0.972	0.903	7.1 (2.1, 11.8)
Airspace Opacity	0.559	0.498	10.9 (-2.0, 23.1)
Edema	0.769	0.746	3.0 (-13.2, 18.5)
Consolidation	0.510	0.738	-44.7 (-130.0, 0.0)
Average	0.814	0.575	29.4 (23.1, 35.5)

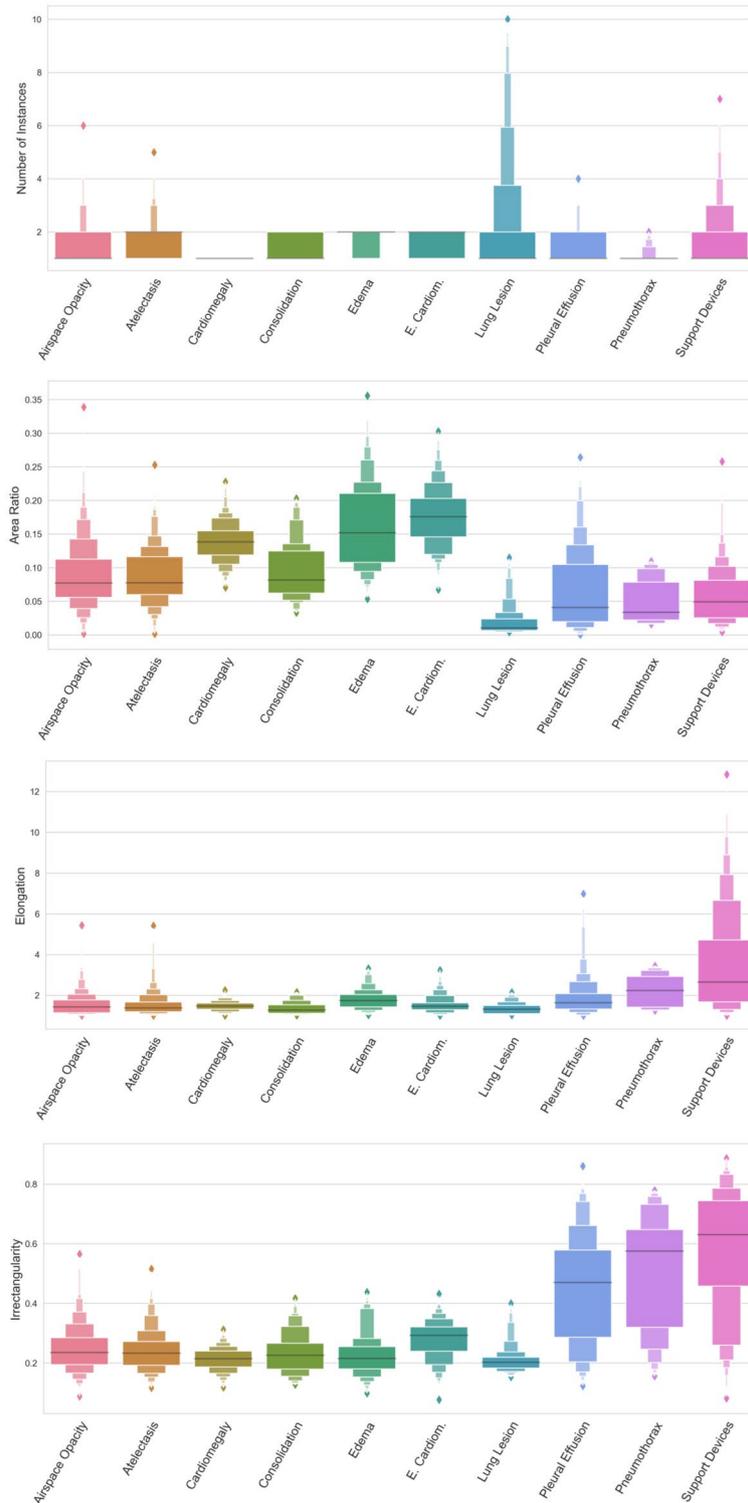
Extended Data Fig. 6 | Percentage decrease from human benchmark hit rate to saliency method pipeline hit rate on the test set. Pathologies are sorted first by statistical significance of percentage decrease from human benchmark

hit rate to saliency method pipeline hit rate (high to low), and then by percentage decrease from human benchmark hit rate to saliency method pipeline hit rate (high to low). We use 95% bootstrap confidence interval.

pathology	ensemble model	single checkpoint
mIoU		
Airspace Opacity	0.248	0.241
Atelectasis	0.254	0.233
Cardiomegaly	0.452	0.419
Consolidation	0.408	0.369
Edema	0.362	0.360
Enlarged Cardiom.	0.379	0.297
Lung Lesion	0.101	0.099
Pleural Effusion	0.235	0.205
Pneumothorax	0.213	0.181
Support Devices	0.163	0.150
hit rate		
Airspace Opacity	0.498	0.534
Atelectasis	0.501	0.504
Cardiomegaly	0.903	0.846
Consolidation	0.738	0.711
Edema	0.746	0.749
Enlarged Cardiom.	0.818	0.704
Lung Lesion	0.290	0.286
Pleural Effusion	0.507	0.390
Pneumothorax	0.392	0.491
Support Devices	0.355	0.312

Extended Data Fig. 7 | Test set saliency method pipeline localization performance using an ensemble model vs. using the top performing single checkpoint for each pathology. For each pathology, we highlight in **bold**

the model (ensemble or single checkpoint) that has the higher metric, and we underline it if the difference is statistically significant (using 95% bootstrap confidence interval).



Extended Data Fig. 8 | Test set distribution of four geometric features across all 10 pathologies. The black horizontal line in each box indicates the median feature value for that pathology, and each successive level outward contains half of the remaining data. The height of the box indicates the range of feature values in the quantile.

pathology	CXRs including all positives and false negatives (<i>n</i>)	Linear regression coefficient	Spearman correlation coefficient
Airspace Opacity	309	1.080 (0.707, 1.453) ***	0.284 (0.178, 0.383) ***
Atelectasis	177	0.481 (-0.086, 1.047)	0.076 (-0.072, 0.221)
Cardiomegaly	175	0.206 (0.006, 0.406) *	0.185 (0.038, 0.324) *
Consolidation	35	1.219 (-0.204, 2.642)	0.324 (-0.011, 0.593)
Edema	83	0.259 (-0.226, 0.745)	0.168 (-0.050, 0.370)
Enlarged Cardiom.	297	0.667 (-0.153, 1.487)	0.191 (0.078, 0.298) ***
Lung Lesion	14	0.657 (-0.825, 2.139)	0.392 (-0.175, 0.764)
Pleural Effusion	120	1.009 (0.533, 1.485) ***	0.347 (0.179, 0.496) ***
Pneumothorax	10	0.342 (-1.549, 2.233)	0.142 (-0.535, 0.708)
Support Devices	314	0.119 (-0.120, 0.358)	0.102 (-0.009, 0.210)
All pathologies	1534	-0.319 (-0.399, -0.239) ***	-0.186 (-0.234, -0.138) ***

* p-value < 0.05, ** p-value < 0.01, *** p-value < 0.001

Extended Data Fig. 9 | Hit/miss: Coefficients from regressions on model assurance. Statistical analysis to determine whether there was any correlation between the model's confidence in its prediction and saliency method pipeline performance using hit/miss. We used the true positive slice of the dataset (CXRs with both the most representative point identified by the saliency method/

human benchmark and also the ground-truth segmentation). Since every heat map contains a maximally activated point (the pixel with the highest value) regardless of model probability output, using the full dataset has limited value since false positives are due to metric set up and are not associated with model probability.

sample size	
Number studies	500
Number CXRs	668
pathology	CXRs (n)
Airspace Opacity	309
Atelectasis	177
Cardiomegaly	175
Consolidation	35
Edema	83
Enlarged Cardiom.	297
Lung Lesion	14
Pleural Effusion	120
Pneumothorax	10
Support Devices	314
No pathology identified	169

Extended Data Fig. 10 | Test set summary statistics.

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- A description of all covariates tested
- A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

Data analysis

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The CheXLocalize dataset is available here: <https://stanforddaimi.azurewebsites.net/datasets/abfb76e5-70d5-4315-badc-c94dd82e3d6d>. The CheXpert dataset is available here <https://stanforddmlgroup.github.io/competitions/chexpert/>.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences Behavioural & social sciences Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://www.nature.com/documents/nr-reporting-summary-flat.pdf)

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	We use 224,314 chest X-ray images from 65,240 patients for model training. For the test set, we selected samples such that we got roughly 20 examples per class.
Data exclusions	No data was excluded.
Replication	The code used to generate segmentations from saliency method heat maps, fine-tune segmentation thresholds, generate segmentations from human annotations, and evaluate localization performance is available in the following public repository under the MIT License: https://github.com/rajpurkarlab/cheXlocalize . The version used for this publication is available at https://doi.org/10.5281/zenodo.681628869 .
Randomization	We didn't require randomization as no human subject evaluation was performed.
Blinding	We didn't require blinding as no human subject evaluation was performed.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	Involvement in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

Methods

n/a	Involvement in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging